

(11)特許出願公開番号

(43)公開日 平成7年(1995)8月18日

H O 4 N 7/15

審査請求 未請求 請求項の数9 OL (全 6 頁)

[最終頁に続く](#)

1

## 【特許請求の範囲】

【請求項1】映像入力手段と、映像表示手段と、映像符号化伝送手段とを含む映像通信装置をネットワークで接続して映像通信を行う装置において、音声入力手段と、前記音声入力手段により入力した音声信号を解析する手段と、前記音声信号の解析手段によって制御される映像符号化手段を備えたことを特徴とする映像通信装置。

【請求項2】映像入力手段と、映像表示手段と、映像符号化伝送手段とを含む映像通信装置をネットワークに接続して映像通信を行う装置において、画像解析手段と、前記画像解析手段により解析した画像信号の特徴に応じて符号生成量が制御される画像符号化手段とを備えたことを特徴とする映像通信装置。

【請求項3】映像入力手段と、映像表示手段と、映像符号化伝送手段とを含む映像通信装置をネットワークに接続して映像通信を行う装置において、音声解析手段と、画像解析手段と、前記音声解析手段と前記画像解析手段とによって符号生成量を制御される画像符号化手段とを備えたことを特徴とする映像通信装置。

【請求項4】請求項1において、前記音声信号を解析する手段が認識する音声信号は、「もしもし」、「おい」、「すみません」、「誰かいませんか」、「応答して下さい」、「やあ」、「さよなら」、「ばいばい」、「ではまた」、「じゃー」からなる呼掛けに使用する言葉である映像通信装置。

【請求項5】請求項1において、前記音声信号を解析する手段が認識する音声信号は、前記映像通信装置の使用人の姓、もしくは名、もしくは姓名、もしくは前記姓名にさん付けしたもの、もしくは役職名を付加したものである映像通信装置。

【請求項6】請求項1において、前記映像信号を解析する手段が解析する画像信号が人間の動作である映像通信装置。

【請求項7】請求項6において、前記人間の動作が「手を振る」、「指を指す」、「画面を見つめる」、「背を向ける」である映像通信装置。

【請求項8】請求項1において、解析すべき前記音声信号を登録可能な映像通信装置。

【請求項9】請求項2において、解析すべき前記映像信号の特徴を登録可能な映像通信装置。

## 【発明の詳細な説明】

## 【0001】

【産業上の利用分野】本発明は、映像によりサテライトオフィス間を接続する映像通信装置に関する。

## 【0002】

【従来の技術】従来より映像通信を利用して会議を行うテレビ会議装置が実用に供している。テレビ会議装置では、カメラで撮影した映像信号をデジタル化し、さらにこの信号を圧縮して伝送している。このように画像を圧縮する理由は以下の通りである。すなわち、テレビ解像

2

度の映像を完全な動画で伝送するには、100メガビット/秒を越える高速な映像伝送ネットワークを使用する必要がある、このような高速回線の伝送コストは極めて高価なためである。そこで、画像信号の持つ統計的特徴を利用して冗長性を取り除く等の手法により100メガビット/秒以上の映像データを、数十キロビット/秒から数メガビット/秒に圧縮して伝送する手法が開発されてきた。現在では、NTTのサービスするISDN回線(Integrated Services on Digital Network: 64キロビット回線)に収容できるほど高効率に圧縮できる技術が開発されてきている。しかし、今までに開発されてきた画像圧縮手法といえども、回線容量内に符号量発生量を抑えるため、駒落しや解像度低下が生じ、臨場感ある映像通信は必ずしも可能とは言えない。また、最近では、さらに伝送容量を増したBISDN(Broadband ISDN)と呼ばれる通信サービスも始まる予定ではあるが、この回線を使用するとしても依然として回線使用量は高く、映像通信サービスを普及させる障害となっている。

## 【0003】

【発明が解決しようとする課題】ところで、首都圏の地価高騰を背景に事務所やオフィスを郊外に設置するような形態(以後、サテライトオフィスと称する)が出現してきている。このサテライトオフィスでは、電話、ファックス、またはコンピュータネットワークを使用してオフィス間を接続し、物理的な距離を感じさせずに作業が可能な環境を提供しようとするもので、1993年時点で幾つかの試行がある。これらの試行を通じてサテライトオフィスをより実用的なものとするには、上記に示した通信装置だけでなく、お互いのオフィスの映像を常時表示して同一空間の共有感を創造する映像通信装置を利用することが効果的であると指摘されている。この理由は、相手地点の映像を常時表示することで同じ空間を共有しているといった感覚を生じさせ、これにより遠隔地のオフィスにいるといった疎外感を軽減することの効果と考えられている。

【0004】ところが、上述のような常時接続する映像通信に、これまでに普及しているテレビ会議装置や映像伝送装置、さらにこれらを接続する回線としてISDN回線やBISDN回線を使用すると、一定量の符号を常に発生するため通信回線を常に占有することになる。従って通信コストが高額になるといった問題がある。

## 【0005】

【課題を解決するための手段】そこで本発明では、通信コストの上昇を抑制しつつ作業効率の向上が図れる臨場感映像通信装置を提供するため、会話の開始と終了を検出するための音声信号解析手段と、画像符号化特性を動画と準動画に切替えられる映像符号化手段と、画像符号化方式に応じて接続する回線を選択する回線接続手段とを設ける。

## 【0006】

3

4

【作用】上記に示した画像符号化手段は、符号発生量は多いが高精細で動きの滑らかな動画信号を発生するモードと、動きに対しては追従性は劣るが符号発生量の少ない符号化モードを備える。音声解析手段は、ユーザの発する音声信号の内容を理解して、必要に応じて符号化特性を切替える。例えば、映像表示画面の前に位置するユーザが臨場感のある会話を行う場合、すなわち、高精細動画モードを必要とする場合は、そのユーザが発生する音声として、「もしもし」、「おーい」などといった言葉を検出し、動画対応の符号化を行うように符号化装置を制御する。この認識装置は、同時に使用する回線を選択し、完全動画モードでは発生する大量のデータを伝送可能な回線に接続するように動作する。一方、会話を終了する際は、会話終了のキーとなる語を検出し、符号化特性を準動画（符号発生量の少ないモード）となるように制御する。これにより、会話時は動きの滑らかで臨場感の高い映像通信が可能となり、また会話のない時は低コストの映像通信が行える。

【0007】

【実施例】図1は本発明による映像通信装置の構成を示したブロック図である。本装置は、音声入力装置1、撮像装置2、画像表示装置3、音声解析装置4、音声符号化装置5、画像符号化装置6、マルチプレクサ9、10、および回線接続装置11から構成される。この映像通信装置は、通信ネットワークを介して遠方にある他の通信装置と接続されている。なお本発明の装置は大容量と少容量のデータを発生する画像符号化装置を備えているため、図1の実施例ではそれぞれの容量に対応した二系統のネットワークを使用することを前提に図を描いている。もちろん回線容量が可変なネットワークでは、一つの回線に異なるデータ量をもつ信号を接続できるので、二つの異なるネットワークに接続する必要がないことは容易に理解できる。なお、この二系統のデータ発生方法については後述する。

【0008】次に図1の装置の動作を説明する。図1の映像通信装置は、撮像装置2で得た映像信号を画像符号化装置6により符号化・圧縮する。映像と同時に収録された音声信号は音声符号化装置5により符号化し、マルチプレクサ9、10により映像信号と合成する。音声信号は、同時に音声解析装置4にも入力しておく。なおこの音声解析装置の動作についても後述する。

【0009】回線接続装置11はこの映像通信装置をネットワークに接続するもので、先に述べた映像・音声信号はネットワークを介して相手の映像通信装置へと送出される。図1の実施例が、従来の映像通信装置と異なる点は、画像符号化装置が二種類の異なるデータレート of の画像符号化装置を有し、発生するデータ量を状況に応じて制御する点である。データレートの異なる画像符号化装置として、ここでは動画符号化装置7と準動画符号化装置8を含む構成として記してある。すなわち、動画符

号化装置7は、駒落しすることなく動きの滑らかな映像を再生するための符号化装置である。従って動画符号化装置は、毎秒に符号化する映像フレーム数が多くなる。そのため、必然的に生成データ量は多くなる。一方、準動画符号化装置8は、駒落しを行うことでデータ発生量の少ない画像符号化を行う装置である。この二系統の符号化装置のいずれを使用するかは制御は音声解析装置4が行う。この音声解析装置4は、同時に回線接続装置11も制御し、この二種類の装置を制御することで、回線に送出する映像・音声データの量を制御する。

【0010】次に、この音声解析装置の動作を詳細に説明する。音声解析装置4は、音声入力装置1からの音声信号を解析し、会話の開始／終了を検出する。会話を開始するなど、臨場感の高い動画符号化に対する要求がある場合は、動画画像符号化装置7により符号化を行い、また要求が無い場合は準動画による符号化を行うように画像符号化装置6を制御する。またこの時、同時に回線接続装置11をも制御し、発生符号量に応じた回線に接続するように動作する。

【0011】図2は音声解析装置4の構成を詳細に示したブロック図である。本音声解析装置は音声照合装置21と登録語データベース22から構成する。この登録語データベース22には動画モードを開始するためのキーワードを登録した動画開始予約語データベース22（図中には動画開始予約語と記載）と動画モードを終了するためのキーワードを登録した動画終了予約語データベース23（図中には動画終了予約語と記載）を備えている。

【0012】信号入力端子20に入力された音声信号は音声照合装置21に入力する。この音声信号は音声認識によりその言葉を解釈する。この言葉（以後、語と記す）は、動画開始予約語データベース23の予約語、動画終了予約語データベースの予約語と照合する。

【0013】動画開始の予約語は、例えば「もしもし」、「おーい」といった呼掛けの言葉や、その他、システムを使用する人名や、その人名に「さん」や肩書を付加したものである。すなわち、会話を開始するために通常良く使用する言葉である。このような語を動画モードの開始予約語として登録しておき、この語を検出した時点で会話状態フラグを会話に設定し、動画による符号化を開始させる。

【0014】一方、動画終了予約語は、例えば、「さよなら」、「じゃー」、「では」、「ばいばい」といった通常会話を終りにする際に使用する語である。このような会話を終了する時に使用する語を登録しておき、この語により会話終了を検出することで動画符号化モードを終了して準動画モードへと移行するようなフラグ信号が発生する。

【0015】このように会話の開始語と、会話の終了語を用いて符号化モードを切替えることの効果を次に説明

する。通常、オフィス映像で結んだ際には、会話が発生していない時に動画を表示することと、駒数を減じた準動画を表示することの差は少なく、動画を伝送しておくことは回線コストがかかるが得られる効果は少ない。従って、会話の無い時は準動画の伝送で十分である。これに対し、通信相手と会話を行おうとする時に駒数の少ない準動画では、動きがぎこちなくなり自然な会話が行えない。そこでこの時は、駒数の豊富な動画信号を使用するのが好ましい。

【0016】音声解析により得られる効果は、会話の開始語と終了語の間にある期間を会話期間と解釈し、この間は高画質な動画符号化により臨場感ある会話環境をサポートし、それ以外の期間は低コストな映像伝送に切替える制御を可能とする。

【0017】図3は画像符号化装置の変形例を表す図である。図1の実施例では、動画符号化装置7と準動画符号化装置8の二系統の符号化装置を用い、これらを音声解析の結果に応じて切替えて使用する構成であったが、図3の変形例では、可変符号化装置により動的に符号発生量を制御するような構成にしてある。この構成によれば、符号化装置を二系統持つ必要がなくなり、回路規模の削減が図れる。この可変符号化装置は、画像符号化装置31、符号化レート制御装置32、レート監視装置33から構成する。レート監視装置33は、音声解析装置4の解析結果に応じて、発生する符号量を決定する。先にも述べたように音声解析装置4は符号生成量を制御するが、それには、レート監視装置33に対して会話開始/終了を指示する会話状態フラグ入力して制御する。

【0018】図4は、図1の実施例の変形を表す図である。図1の実施例と図4の実施例の大きく異なる点は画像解析装置40を設けた点である。この画像解析装置40によって画像符号化装置43、および回線接続装置41を制御する点が新規な点である。図3の画像解析装置40を設けた目的は、音声認識装置4を用いるのと同様であり、会話の開始/終了を検出するためである。そのため、画像解析装置は撮像した映像の中から、会話の開始を意味する動作（以下、ジェスチャと記す）を検出し、このジェスチャに基づいて画像符号化装置43と回線接続装置41を制御する。制御方法は先ほどの音声解析の場合と同様に、会話期間は高画質の動画を伝送するように、また会話の終了以後は符号量の少ない準動画を伝送するように制御する。

【0019】図5は画像解析装置40の構成を示した実施例である。画像解析装置40は画像照合装置52とジェスチャデータベース52から構成する。このジェスチャデータベース52は動画開始予約データベース53と動画終了データベース54からなる。この動画開始予約データベースには、会話の開始を意味するジェスチャとして、「手を振る」、「画面を見つめる」、「指で指す」、「手招きをする」などの動作を登録しておく。こ

のように登録したジェスチャにより会話の開始を検出し、会話が始まった場合には会話の開始を表すフラグ信号を発生する。また、動画終了予約ジェスチャには、「手を振る」、「背を向ける」などのジェスチャを登録しておき、このようなジェスチャを認識した時に会話の終了を表す終了フラグを発生する。

【0020】なお、ジェスチャの照合方式は本発明の趣旨とは直接関係ないので、ここでの詳細な説明は省略するが、例えば「手を振る」動作を認識するには、手の形のテンプレートをデータベースに用意しておき、このテンプレートとのパターンマッチングにより手を認識し、さらにこの手の動きベクトルを検査することで、「手を振る」といった動作が解析できる。なお、この説明はジェスチャ認識のほんの一例であり、どのような認識方法であっても本発明の趣旨が満足されることは容易に理解できる。

【0021】また、図5の実施例では、会話の開始と終了と共に「手を振る」動作が含まれているが、これはトグル操作の意味であり、始めに「手を振る」動作が会話の開始、次の「手を振る」動作が会話の終了といったように解釈すればよい。

【0022】さらに、本装置による使い勝手を向上させるために、図2の音声解析装置のキーワードや、図5の画像解析装置のキーとなるジェスチャはユーザが登録可能な構成にしておくことが好ましい。すなわち、通常良く使用する言葉やジェスチャを必要に応じて登録することで、会話の開始/終了の照合をより確実とすることが可能となる。

【0023】本実施例によれば音声信号、あるいは表示画面の前に立った人物のジェスチャを解析することで会話の始めと終りを検出し、会話が行われている期間だけ画像符号化装置を高画質な動画モードで動作させる。一方、会話の行われていない期間は、駒数を減じてデータ発生量の少ない準動画符号化とする。

【0024】

【発明の効果】本発明によれば、ユーザが画面に写った相手オフィスの映像に呼掛けを行なった時だけコストのかかる伝送路を使用し、その他の場合には通信コストの低い回線を使用する。その結果として会話の時の臨場感を確保しつつ、通信コストの少ない映像通信が可能となる。

【図面の簡単な説明】

【図1】本発明による映像通信装置のブロック図。

【図2】音声解析装置のブロック図。

【図3】映像符号化装置の変形例のブロック図。

【図4】ジェスチャ解析装置を含む映像通信装置のブロック図。

【図5】ジェスチャ解析装置のブロック図。

【符号の説明】

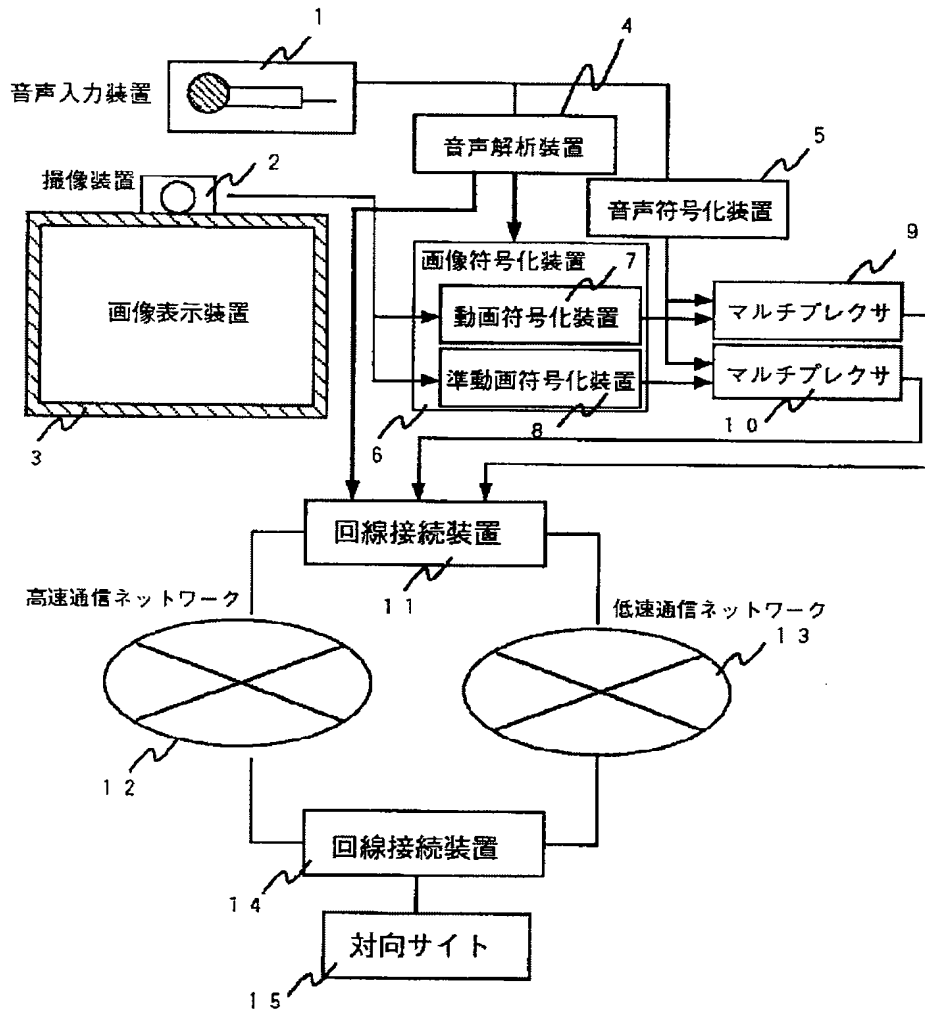
1…音声入力装置、2…撮像装置、3…表示装置、4…

音声解析装置、5…音声符号化装置、6…画像符号化装置、7…動画符号化装置、8…準動画符号化装置、9、10…マルチプレкса、11、14…回線接続装置、1

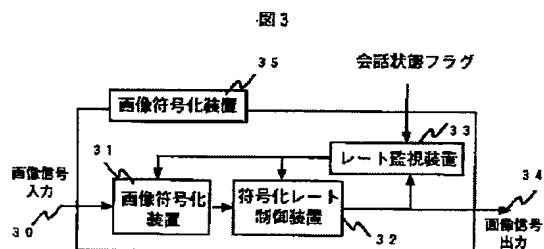
2…高速ネットワーク、13…低速ネットワーク、15…対向サイト。

【図1】

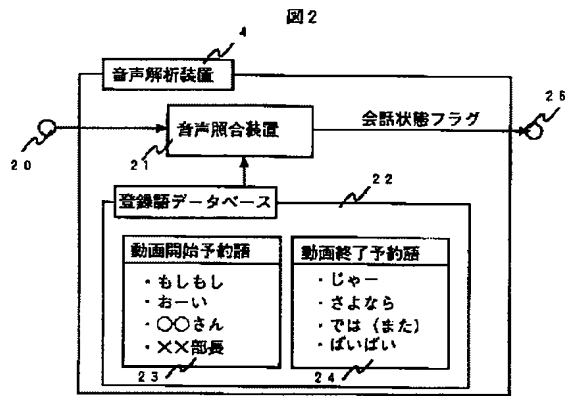
図1



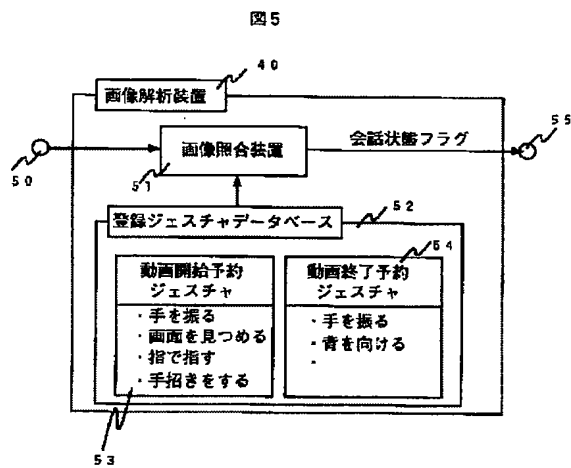
【図3】



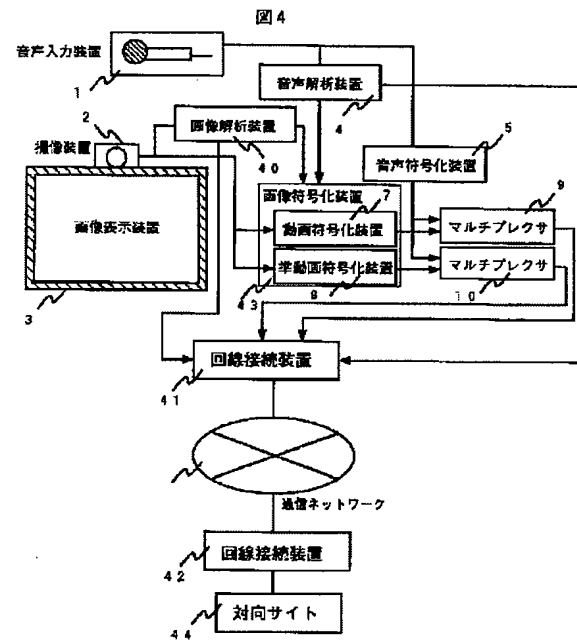
【図2】



【図5】



【図4】



フロントページの続き

(72)発明者 木下 泰三

東京都国分寺市東恋ヶ窪1丁目280番地  
株式会社日立製作所中央研究所内